

Analysis of single-cell RNA sequencing data based on autoencoders

Andrea Tangherloni, Federico Ricciuti, Daniela Besozzi,
Pietro Liò, Ana Cvejic

Mathematical formulation of the proposed autoencoders

In what follows, we use the notation proposed in [1] to derive the extension of the Mean Maximum Discrepancy Autoencoder (MMDAE) and Mean Maximum Discrepancy Variational Autoencoder (MMDVAE) [1] with multiple Gaussian distributions.

$p^*(\mathbf{x})$ is the unknown probability in the input space over which the optimisation problem is formulated. A similarity measure \mathcal{L} between the input and the output of the autoencoders (AEs) with respect to the distribution $p^*(\mathbf{x})$ is maximised:

$$\arg \max_{\phi, \theta} \mathbb{E}[\mathcal{L}(\mathbf{x}, d_{\theta}(e_{\phi}(\mathbf{x})))] ,$$

where ϕ and θ are the weights of the encoder and decoder networks, respectively; $e_{\phi} : \mathbf{x} \mapsto \mathbf{z}$ and $d_{\theta} : \mathbf{z} \mapsto \mathbf{x}$, where \mathbf{z} is the latent representation of \mathbf{x} and $|\mathbf{z}| \leq |\mathbf{x}|$.

In variational AEs, the input \mathbf{x} is mapped into a probability distribution over the latent space. $e(\mathbf{z}|\mathbf{x})$ defines a distribution over the latent space that depends on the input \mathbf{x} drawn from $p^*(\mathbf{x})$. Altogether, $p^*(\mathbf{x})$ and $e(\mathbf{z}|\mathbf{x})$ define the joint distribution $p_e(\mathbf{x}, \mathbf{z}) = e(\mathbf{z}|\mathbf{x})p^*(\mathbf{x})$, whose marginal and conditional distributions are defined as:

$$p_e(\mathbf{z}) = \int p_e(\mathbf{x}, \mathbf{z}) d\mathbf{x} = \int p_e(\mathbf{z}|\mathbf{x}) p^*(\mathbf{x}) d\mathbf{x}$$
$$p_e(\mathbf{x}|\mathbf{z}) = \frac{p_e(\mathbf{x}, \mathbf{z})}{p_e(\mathbf{z})}.$$

Since the representation \mathbf{z} of \mathbf{x} should maintain as much as possible the “amount of information” held in \mathbf{x} , the mutual information $I(\mathbf{x}; \mathbf{z})$ can be used to measure the representation \mathbf{z} of \mathbf{x} . Specifically, for any distribution $q(\mathbf{z})$ in the latent space, the mutual information between $p_e(\mathbf{z})$ and $p^*(\mathbf{x})$ can be bounded below as:

$$I(\mathbf{x}; \mathbf{z}) = \text{KL}(p_e(\mathbf{x}, \mathbf{z}) || p^*(\mathbf{x}) p_e(\mathbf{z})) \leq \mathbb{E}[\text{KL}(e(\mathbf{z}|\mathbf{x}) || q(\mathbf{z}))],$$

where $\text{KL}(\cdot)$ is the Kullback–Leibler divergence [2] between two distributions. $I(\mathbf{x}; \mathbf{z})$ can be also bounded above, for any conditional distribution $d(\mathbf{x}|\mathbf{z})$, as:

$$I(\mathbf{x}; \mathbf{z}) = \text{KL}(p_e(\mathbf{x}, \mathbf{z}) || p^*(\mathbf{x}) p_e(\mathbf{z})) \geq \mathbb{E} \left[\log \left(\frac{d(\mathbf{x}|\mathbf{z})}{p^*(\mathbf{x})} \right) \right].$$

Combining the provided definitions, we obtain that:

$$\mathbb{E} \left[\log \left(\frac{d(\mathbf{x}|\mathbf{z})}{p^*(\mathbf{x})} \right) \right] \leq I(\mathbf{x}; \mathbf{z}) \leq \mathbb{E}[\text{KL}(e(\mathbf{z}|\mathbf{x})||q(\mathbf{z}))].$$

The lower bound can be further decomposed by means of algebraic manipulations as

$$\mathbb{E} \left[\log \left(\frac{d(\mathbf{x}|\mathbf{z})}{p^*(\mathbf{x})} \right) \right] = \mathbb{E}[\log(d(\mathbf{x}|\mathbf{z})) + H(p^*(\mathbf{x}))],$$

where $H(p^*(\mathbf{x}))$ is the entropy of $p^*(\mathbf{x})$. By following the definition provided in [1], the ELBO term, which is the measure maximised during the training of VAE, can be written as:

$$\text{ELBO} = -\text{KL}(p_e(\mathbf{z})||q(\mathbf{z})) - H(p^*(\mathbf{x})) - \mathbb{E}[\text{KL}(p_e(\mathbf{x}|\mathbf{z})||d(\mathbf{x}|\mathbf{z}))].$$

In MMDVAE [1], $\text{KL}(p_e(\mathbf{z})||q(\mathbf{z}))$ is multiplied by a positive factor λ and $I(\mathbf{x}; \mathbf{z})$, weighted by a positive factor α , is added to the ELBO term, obtaining:

$$\begin{aligned} \text{ELBO} = & -\lambda \text{KL}(p_e(\mathbf{z})||q(\mathbf{z})) \\ & - H(p^*(\mathbf{x})) \\ & - \mathbb{E}[\text{KL}(p_e(\mathbf{x}|\mathbf{z})||d(\mathbf{x}|\mathbf{z}))] \\ & + \alpha I(\mathbf{x}; \mathbf{z}). \end{aligned}$$

By applying algebraic manipulations, the ELBO term of MMDVAE can be written as:

$$\begin{aligned} \text{ELBO} = & \mathbb{E}[\log(d(\mathbf{x}|\mathbf{z}))] \\ & - (\alpha + \lambda - 1)\text{KL}(p_e(\mathbf{z})||q(\mathbf{z})) \\ & - (1 - \alpha)\mathbb{E}[\text{KL}(p_e(\mathbf{z}|\mathbf{x})||q(\mathbf{z}))]. \end{aligned} \tag{1}$$

In MMDVAE, the term $\text{KL}(p_e(\mathbf{z})||q(\mathbf{z}))$ is replaced with $\text{DSD}(p_e(\mathbf{z})||q(\mathbf{z}))$, where $\text{DSD}(\cdot)$ is a general strict divergence function. $\text{DSD}(p_e(\mathbf{z})||q(\mathbf{z})) = 0$ if and only if $p_e(\cdot) = q(\cdot)$. Notice that, the KL is a strict divergence function. MMDVAE exploits the Maximum Mean Discrepancy $\text{MMD}(\cdot)$ divergence function [3]. A kernel trick is used to define the following divergence function between two distributions $p_e(\mathbf{z})$ and $q(\mathbf{z})$:

$$\begin{aligned} \text{MMD}(p_e(\mathbf{z})||q(\mathbf{z})) = & \mathbb{E}_{p_e(\mathbf{z}), p(\mathbf{z}')} [\mathcal{K}(\mathbf{z}, \mathbf{z}')] \\ & + \mathbb{E}_{q(\mathbf{z}), q(\mathbf{z}')} [\mathcal{K}(\mathbf{z}, \mathbf{z}')] \\ & - 2\mathbb{E}_{p_e(\mathbf{z}), q(\mathbf{z}')} [\mathcal{K}(\mathbf{z}, \mathbf{z}')], \end{aligned}$$

where $\mathcal{K}(\mathbf{z}, \mathbf{z}')$ can be any desired universal kernel. Here, we considered the Gaussian kernel

$$\mathcal{K}(\mathbf{z}, \mathbf{z}') = e^{-\frac{\|\mathbf{z}-\mathbf{z}'\|}{2\sigma^2}}.$$

We extended the ELBO term shown in Eq. 1 such that multiple Gaussian distributions can be used in the latent representation \mathbf{z} . In addition, we introduced a learnable mixture distribution for $q(\mathbf{z})$, whereas $p_e(\mathbf{z}|\mathbf{x})$ is defined to be a learnable mixture distribution with the same number of components.

In GMVAE [4], the encoder function outputs the following two conditional distributions $e(\mathbf{z}, y|\mathbf{x})$ and $e(\mathbf{z}|\mathbf{x}, y)$, where $y \in \{1, \dots, K\}$ is a categorical random variable and K corresponds to the number of desired Gaussian distributions. We obtain that

$$\begin{aligned} p_e(\mathbf{z}, y|\mathbf{x}) &= \frac{p_e(\mathbf{z}, y, \mathbf{x})}{p^*(\mathbf{x})} \\ &= e(\mathbf{z}|y, \mathbf{x})e(y|\mathbf{x}). \end{aligned}$$

Since $p_e(\mathbf{z}, y|\mathbf{x})$ is fully determined by the output distributions of the encoder, we can refer to $p_e(\mathbf{z}, y|\mathbf{x})$ with $e(\mathbf{z}, y|\mathbf{x})$.

Modelling $e(y|\mathbf{x})$ as a categorical distribution that can assume values in $\{1, \dots, K\}$, and $e(\mathbf{z}|\mathbf{x}, y)$ as a diagonal Gaussian distribution for each possible value assumed by y , the marginal conditional distribution $p_e(\mathbf{z}|\mathbf{x})$ is a Gaussian mixture distribution of K components, namely:

$$p_e(\mathbf{z}|\mathbf{x}) = \sum_{y=1}^K e(\mathbf{z}|y, \mathbf{x})e(y|\mathbf{x}).$$

Similarly, $q(\mathbf{z})$ is modelled as a Gaussian mixture distribution by using another variable $y \in \{1, \dots, K\}$, with a categorical distribution $q(y)$, and considering the conditional distribution $q(\mathbf{z}|y)$ as a diagonal Gaussian distribution for each possible value of y . The ELBO term of GMVAE is:

$$ELBO = \mathbb{E} \left[\mathbb{E} \left[d(\mathbf{x}|y, \mathbf{z}) - \log \left(\frac{e(\mathbf{z}, y|\mathbf{x})}{q(\mathbf{z}, y)} \right) \right] \right],$$

and it can be rewritten by using the notation proposed in [1] and by algebraic manipulations as:

$$\begin{aligned} ELBO &= -\text{KL}(p_e(\mathbf{z}, y)||q(\mathbf{z}, y)) \\ &= -\mathbb{E}[\text{KL}(p_e(\mathbf{x}|\mathbf{z}, y)||d(\mathbf{x}|\mathbf{z}, y))] \\ &= -H(p^*(\mathbf{x})). \end{aligned}$$

Starting from this definition, we can add the mutual information $I(\mathbf{x}; (y, \mathbf{z}))$ term, weighted by a positive scalar factor α , and $\text{KL}(p_e(\mathbf{z}, y)||q(\mathbf{z}, y))$ is weighted by a positive factor λ , obtaining:

$$\begin{aligned} ELBO &= -\lambda \text{KL}(p_e(\mathbf{z}, y)||q(\mathbf{z}, y)) \\ &\quad - H(p^*(\mathbf{x})) \\ &\quad - \mathbb{E}[\text{KL}(p_e(\mathbf{x}|\mathbf{z}, y)||d(\mathbf{x}|\mathbf{z}, y))] \\ &\quad + \alpha I(\mathbf{x}; (y, \mathbf{z})), \end{aligned}$$

where

$$I(\mathbf{x}; (y, \mathbf{z})) = \mathbb{E} \left[\log \frac{p_e(\mathbf{x}, y, \mathbf{z})}{p^*(\mathbf{x})p_e(\mathbf{z}, y)} \right].$$

By applying algebraic manipulations, we can rewrite the ELBO term as:

$$\begin{aligned} ELBO &= \mathbb{E}[\log(d(\mathbf{x}|\mathbf{z}, y))] \\ &\quad - (\alpha + \lambda - 1) \text{KL}(p_e(\mathbf{z}, y)||q(\mathbf{z}, y)) \\ &\quad - (1 - \alpha) \mathbb{E}[\text{KL}(p_e(\mathbf{z}, y|\mathbf{x})||q(\mathbf{z}, y))]. \end{aligned}$$

$\text{KL}(p_e(\mathbf{z}, y) || q(\mathbf{z}, y))$ can be further decomposed as:

$$\text{KL}(p_e(\mathbf{z}, y) || q(\mathbf{z}, y)) = \mathbb{E}[\text{KL}(p_e(y|\mathbf{z}) || q(y|\mathbf{z}))] + \text{KL}(p_e(\mathbf{z}) || q(\mathbf{z})),$$

so that the ELBO can be written as:

$$\begin{aligned} \text{ELBO} &= \mathbb{E}[\log(d(\mathbf{x}|\mathbf{z}, y))] \\ &\quad - (\alpha + \lambda - 1) \text{KL}(p_e(\mathbf{z}) || q(\mathbf{z})) \\ &\quad - (1 - \alpha) \mathbb{E}[\text{KL}(p_e(\mathbf{z}, y|\mathbf{x}) || q(\mathbf{z}, y))]. \end{aligned}$$

As in MMDVAE (see Eq. 1), we can replace $\text{KL}(p_e(\mathbf{z}) || q(\mathbf{z}))$ with a general strict divergence function. We considered the $\text{MMD}(\cdot)$ term, obtaining the a general formulation for all the five AEs:

$$\begin{aligned} \text{ELBO} &= \mathbb{E}[\log(d(\mathbf{x}|\mathbf{z}, y))] \\ &\quad - (\alpha + \lambda - 1) \text{MMD}(p_e(\mathbf{z}) || q(\mathbf{z})) \\ &\quad - (1 - \alpha) \mathbb{E}[\text{KL}(p_e(\mathbf{z}, y|\mathbf{x}) || q(\mathbf{z}, y))]. \end{aligned} \tag{2}$$

We modified the $\text{MMD}(p_e(\mathbf{z}) || q(\mathbf{z}))$ such that it is not necessary to sample from the Gaussian mixture distribution $e(\mathbf{z}|\mathbf{x})$ or from the posterior $q(\mathbf{z})$. Our modification allows for sampling from the single Gaussian distributions that form the mixtures. We used the the reparametrization trick proposed in [5] so that $\text{MMD}(p_e(\mathbf{z}) || q(\mathbf{z}))$ can be approximated. Specifically, $\mathbb{E}_{p_e(\mathbf{z}), p(\mathbf{z}')}[\mathcal{K}(\mathbf{z}, \mathbf{z}')] can be approximated as:$

$$\frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N \sum_{y=1}^K \sum_{y'=1}^K p_e(y, x_i) p_e(y', x'_j) \mathbb{E}_{p_e(\mathbf{z}|x_i, y), p_e(\mathbf{z}'|x'_j, y')}[\mathcal{K}(\mathbf{z}, \mathbf{z}')],$$

where N is the number of provided samples (i.e., cells). The approximation of $\mathbb{E}_{q(\mathbf{z}), q(\mathbf{z}')}[\mathcal{K}(\mathbf{z}, \mathbf{z}')] is:$

$$\sum_{y=1}^K \sum_{y'=1}^K q(y) q(y') \mathbb{E}_{q(\mathbf{z}|y), q(\mathbf{z}'|y')}[\mathcal{K}(\mathbf{z}, \mathbf{z}')].$$

Finally, $\mathbb{E}_{p_e(\mathbf{z}), q(\mathbf{z}')}[\mathcal{K}(\mathbf{z}, \mathbf{z}')] is approximated as:$

$$\sum_{i=1}^N \sum_{y=1}^K \sum_{y'=1}^K p_e(y|x_i) q(y') \mathbb{E}_{p_e(\mathbf{z}|y, x_i), q(\mathbf{z}'|y')}[\mathcal{K}(\mathbf{z}, \mathbf{z}')].$$

To calculate the ELBO function described in Eq. 2, $\mathbb{E}[\text{KL}(p_e(\mathbf{z}, y|\mathbf{x}) || q(\mathbf{z}, y))]$ must be computable. We can rewrite it as follows:

$$\begin{aligned} \mathbb{E}[\text{KL}(p_e(\mathbf{z}, y|\mathbf{x}) || q(\mathbf{z}, y))] &= \mathbb{E} \left[\log \frac{p_e(\mathbf{z}|y, \mathbf{x}) p_e(y|\mathbf{x})}{q(\mathbf{z}|y) q(y)} \right] \\ &= \mathbb{E} \left[\log \frac{p_e(\mathbf{z}|y, \mathbf{x})}{q(\mathbf{z}|y)} \right] + \mathbb{E} \left[\mathbb{E} \left[\log \frac{p_e(y|\mathbf{x})}{q(y)} \right] \right] \\ &= \mathbb{E}[\text{KL}(p_e(y|\mathbf{x}) || q(y))] + \mathbb{E}[\mathbb{E}[\text{KL}(p_e(\mathbf{z}|\mathbf{x}, y) || q(\mathbf{z}|y))]]. \end{aligned}$$

Considering the weights of $q(\mathbf{z})$ fixed to a uniform distribution, $\mathbb{E}[\text{KL}(p_e(y|\mathbf{x})||q(y))]$ can be written as:

$$\mathbb{E}[\text{KL}(p_e(y|\mathbf{x})||q(y))] = \mathbb{E}\left[\sum_{y=1}^K p_e(y|\mathbf{x}) \log(p_e(y|\mathbf{x}))\right] + \log(K).$$

On the contrary, when the weights are learnable, $\mathbb{E}[\text{KL}(p_e(y|\mathbf{x})||q(y))]$ can be analytically calculated as:

$$\mathbb{E}[\text{KL}(p_e(y|\mathbf{x})||q(y))] = \mathbb{E}\left[\sum_{y=1}^K p_e(y|\mathbf{x}) \log(p_e(y|\mathbf{x})) - \sum_{y=1}^K p_e(y|\mathbf{x}) \log(q(y))\right].$$

Finally, $\text{KL}(p_e(\mathbf{z}|\mathbf{x}, y)||q(\mathbf{z}|y))$ can be calculated by the following approximation:

$$\text{KL}(p_e(\mathbf{z}|\mathbf{x}, y)||q(\mathbf{z}|y)) = \mathbb{E}\left[\log \frac{p_e(\mathbf{z}|\mathbf{x}, y)}{q(\mathbf{z}|y)}\right] \approx \log \frac{p_e(\mathbf{z}|\mathbf{x}, y)}{q(\mathbf{z}|y)}$$

Additional Files

Additional file 2 — Excel file of the metrics calculated for the PBMC datasets

Each tab is related to a tested approach and shows the calculated metrics and used method.

Additional file 3 — Excel file of the metrics calculated for the PIC datasets

Each tab is related to a tested approach and shows the calculated metrics and used method.

Additional file 4 — Excel file of the metrics calculated for the MCA datasets

Each tab is related to a tested approach and shows the calculated metrics and used method.

References

- [1] Zhao, S., Song, J., Ermon, S.: Infovae: Balancing learning and inference in variational autoencoders. In: Proceedings of the AAAI Conference on Artificial Intelligence, vol. 33, pp. 5885–5892 (2019)
- [2] Kullback, S., Leibler, R.A.: On information and sufficiency. Ann Math Statist. **22**(1), 79–86 (1951). doi:10.1214/aoms/1177729694
- [3] Gretton, A., Borgwardt, K., Rasch, M., Schölkopf, B., Smola, A.J.: A kernel method for the two-sample-problem. In: Proceedings of the Conference on Advances in Neural Information Processing Systems, pp. 513–520 (2007)

- [4] Grønbech, C.H., Vording, M.F., Timshel, P.N., Sønderby, C.K., Pers, T.H., Winther, O.: scVAE: Variational auto-encoders for single-cell gene expression data. *Bioinformatics* (2020). doi:10.1093/bioinformatics/btaa293
- [5] Kingma, D.P., Welling, M.: Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114* (2013)